

the university of edinburgh

Applied Machine Learning (AML)

Non-Linear Dimensionality Reduction

Oisin Mac Aodha • Siddharth N.

Outline

Non-Linearity

- Extensions for Linear Dimensionality Reduction
 - Kernel PCA
- Visualisation
 - Multi-Dimensional Scaling (MDS)
 - Isomap
 - Locally Linear Embeddings (LLE)
 - t-distributed Stochastic Neighbour Embedding (t-SNE)
 - Uniform Manifold Approximation & Projection (UMAP)



Extensions for Linear Dimensionality Reduction

PCA on Non-Linear Data

Example: Shells



PCA on Non-Linear Data

Example: Shells



Key Idea: Transform inputs \boldsymbol{x} using a feature map $\phi(\boldsymbol{x})$



Key Idea: Transform inputs \boldsymbol{x} using a feature map $\phi(\boldsymbol{x})$ $\boldsymbol{x} \in \mathbb{R}^{D}, \ \phi(\boldsymbol{x}) \in \mathbb{R}^{C}, \ C > D!$



Key Idea: Transform inputs x using a feature map $\phi(x)$ $x \in \mathbb{R}^{D}, \ \phi(x) \in \mathbb{R}^{C}, \ C > D!$





Key Idea: Transform inputs \boldsymbol{x} using a feature map $\phi(\boldsymbol{x})$ $\boldsymbol{x} \in \mathbb{R}^{D}, \ \phi(\boldsymbol{x}) \in \mathbb{R}^{C}, \ C > D!$





 $X = [\mathbf{x}_1; \dots; \mathbf{x}_N],$ $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ (kernel function) $\Phi(X) = [\phi(\mathbf{x}_1); \dots; \phi(\mathbf{x}_N)]$ $K = \Phi(X)^\top \Phi(X)$

(kernel matrix)



$$X = [\mathbf{x}_1; \dots; \mathbf{x}_N],$$

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

(kernel function)

$$\Phi(X) = [\phi(\boldsymbol{x}_1); \dots; \phi(\boldsymbol{x}_N)]$$
$$K = \Phi(X)^{\top} \Phi(X)$$

(kernel matrix)

$$S = \frac{1}{N} X X^{\mathsf{T}}$$
$$S \boldsymbol{v} = \lambda \boldsymbol{v}$$



$$X = [\mathbf{x}_1; \dots; \mathbf{x}_N],$$

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

(kernel function)

$$\Phi(X) = [\phi(\mathbf{x}_1); \dots; \phi(\mathbf{x}_N)]$$
$$K = \Phi(X)^{\top} \Phi(X)$$

(kernel matrix)

$$S = \frac{1}{N} X X^{\top}$$
$$S \boldsymbol{v} = \lambda \boldsymbol{v}$$

$$S = \frac{1}{N} \Phi(X) \Phi(X)^{\top}$$
$$\boldsymbol{v} = \Phi(X) \boldsymbol{a} = \sum_{i=1}^{N} a_{i} \phi(\boldsymbol{x}_{i})$$
(...some algebra)

$$K\boldsymbol{a} = N\lambda\boldsymbol{a}$$



$$X = [\mathbf{x}_1; \dots; \mathbf{x}_N],$$
$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$
(kernel function)

$$\Phi(X) = [\phi(\boldsymbol{x}_1); \dots; \phi(\boldsymbol{x}_N)]$$
$$K = \Phi(X)^{\top} \Phi(X)$$

(kernel matrix)

$$S = \frac{1}{N} X X^{\top}$$
$$S \boldsymbol{v} = \lambda \boldsymbol{v}$$

Kernel Trick

No need to compute $\phi(\mathbf{x})$ —only $\kappa(\mathbf{x}_i, \mathbf{x}_j)$!



$$K\boldsymbol{a} = N\lambda \boldsymbol{a}$$



Kernel PCA: Example

$$\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \exp(\frac{-\|\boldsymbol{x}_i - \boldsymbol{x}_j\|^2}{2\sigma^2})$$



Kernel PCA: Example





Choosing the right kernel

• Not an easy choice



Choosing the right kernel

- Not an easy choice
- Construct by hand/eye where feasible



Choosing the right kernel

- Not an easy choice
- Construct by hand/eye where feasible
- Possible to learn kernel matrices *K* directly from data!



Choosing the right kernel

- Not an easy choice
- Construct by hand/eye where feasible
- Possible to learn kernel matrices *K* directly from data!



Choosing the right kernel

- Not an easy choice
- Construct by hand/eye where feasible
- Possible to learn kernel matrices *K* directly from data!

Several non-linear dimensionality reduction methods can be viewed as kernel PCA, with kernels learned from data [1]





Visualisation

Manifold Hypothesis

High-dimensional data in the real world really lies on low-dimensional manifolds within that high-dimensional space.



Manifold Hypothesis

High-dimensional data in the real world really lies on low-dimensional manifolds within that high-dimensional space.





Manifold Hypothesis

High-dimensional data in the real world really lies on low-dimensional manifolds within that high-dimensional space.





Overview

Key Ideas

- Difficult to construct a single global transformation of data
- Focus instead on some *local* measure of closeness
- Project data *x* onto lower-dimensional manifold as *e*
- Question: can we *preserve* local measure of closeness?



Overview

Key Ideas

- Difficult to construct a single global transformation of data
- Focus instead on some *local* measure of closeness
- Project data *x* onto lower-dimensional manifold as *e*
- Question: can we preserve local measure of closeness?

Let X denote the high-dimensional data space, and \mathcal{E} denote the low-dimensional manifold space. We can define the following distance measures on the two spaces respectively

$$\mathcal{D}_{\mathcal{X}}(oldsymbol{x}_i,oldsymbol{x}_j) \qquad \qquad \mathcal{D}_{\mathcal{E}}(oldsymbol{e}_i,oldsymbol{e}_j)$$





$$\mathcal{D}_{\mathcal{X}}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \|\boldsymbol{x}_i - \boldsymbol{x}_j\| \qquad \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_i, \boldsymbol{e}_j) = \|\boldsymbol{e}_i - \boldsymbol{e}_j\| \qquad (example)$$



$$\mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) = \|\boldsymbol{x}_{i} - \boldsymbol{x}_{j}\| \qquad \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_{i}, \boldsymbol{e}_{j}) = \|\boldsymbol{e}_{i} - \boldsymbol{e}_{j}\| \qquad (\text{example})$$

Objective:
$$\min \sum_{i,j} (\mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) - \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_{i}, \boldsymbol{e}_{j}))^{2}$$



$$\mathcal{D}_{\mathcal{X}}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) = \|\boldsymbol{x}_{i} - \boldsymbol{x}_{j}\| \qquad \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_{i}, \boldsymbol{e}_{j}) = \|\boldsymbol{e}_{i} - \boldsymbol{e}_{j}\| \qquad (\text{example})$$

Objective:
$$\min \sum_{i,j} (\mathcal{D}_{\mathcal{X}}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) - \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_{i}, \boldsymbol{e}_{j}))^{2}$$

- distance metrics $\mathcal{D}_{\chi}, \mathcal{D}_{\mathcal{E}}$ could really be anything
- choosing L₂ helps make optimisation simpler



MDS: Example

Swiss Roll





MDS: Example

Swiss Roll



- Projects down to 2D while preserving distances
- Preserving distances outside the manifold!



Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!





Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



Find the Geodesic

• Generate nearest-neighbour graph *G* on data



Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



Find the Geodesic

- Generate nearest-neighbour graph *G* on data
- Shortest distance between points in this graph
 - Floyd-Warshall algorithm (all pairs shortest path)


Isomap

Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



Find the Geodesic

- Generate nearest-neighbour graph *G* on data
- Shortest distance between points in this graph
 - Floyd-Warshall algorithm (all pairs shortest path)

• Perform MDS with
$$\mathcal{D}_{\mathcal{X}}(x_i, x_j) = \mathcal{G}_{\mathsf{FW}}(x_i, x_j)$$



Isomap

Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



Find the Geodesic

- Generate nearest-neighbour graph *G* on data
- Shortest distance between points in this graph
 - Floyd-Warshall algorithm (all pairs shortest path)

• Perform MDS with
$$\mathcal{D}_{\mathcal{X}}(x_i, x_j) = \mathcal{G}_{\mathsf{FW}}(x_i, x_j)$$



Isomap

Project data from X to \mathcal{E} while preserving the distance between points on the *embedded manifold*, not arbitrary distance!



Find the Geodesic

- Generate nearest-neighbour graph *G* on data
- Shortest distance between points in this graph
 - Floyd-Warshall algorithm (all pairs shortest path)

• Perform MDS with
$$\mathcal{D}_{\mathcal{X}}(x_i, x_j) = \mathcal{G}_{\mathsf{FW}}(x_i, x_j)$$

Objective:
$$\min \sum_{i,j} (\mathcal{G}_{\mathsf{FW}}(\pmb{x}_i, \pmb{x}_j) - \|\pmb{e}_i - \pmb{e}_j\|)^2$$



Isomap: Manifold



Data



Isomap: Manifold



+ Floyd-Warshall



Isomap: Manifold



+ Floyd-Warshall







the university of edinburgh



Issues

• requires *uniform* and *dense* sampling on manifold



Issues

- requires *uniform* and *dense* sampling on manifold
- prone to topological instabilities (effect of noise, non-convexity)





Issues

- requires *uniform* and *dense* sampling on manifold
- prone to topological instabilities (effect of noise, non-convexity)
- can get disconnected graphs!





Issues

- requires *uniform* and *dense* sampling on manifold
- prone to topological instabilities (effect of noise, non-convexity)
- can get disconnected graphs!
- slow with size of data —
 Floyd-Warshall is O(N³)!



Locally Linear Embeddings (LLE)

Project data from X to \mathcal{E} while preserving the *linear transform that reconstructs points* from the *K* nearest neighbours



Locally Linear Embeddings (LLE)

Project data from X to \mathcal{E} while preserving the *linear transform that reconstructs points* from the *K* nearest neighbours

$$\arg \min_{w_1,...,w_K} \sum_{i=1}^N (\boldsymbol{x}_i - \sum_{\boldsymbol{x}_k \in n(\boldsymbol{x}_i)} w_k \boldsymbol{x}_k)^2$$
$$\arg \min_{\boldsymbol{e}_1,...,\boldsymbol{e}_N} \sum_{i=1}^N (\boldsymbol{e}_i - \sum_{\boldsymbol{e}_k \in n(\boldsymbol{e}_i)} w_k \boldsymbol{e}_k)^2$$

Algorithm

- **1**. Find the weights w_k
- 2. Fix weights w_k and find optimal embeddings e_i solved using a (sparse) eigen-decomposition





Locally Linear Embeddings (LLE)

Project data from X to \mathcal{E} while preserving the *linear transform that reconstructs points* from the *K* nearest neighbours

$$\arg \min_{w_1,...,w_K} \sum_{i=1}^N (\boldsymbol{x}_i - \sum_{\boldsymbol{x}_k \in n(\boldsymbol{x}_i)} w_k \boldsymbol{x}_k)^2$$
$$\arg \min_{\boldsymbol{e}_1,...,\boldsymbol{e}_N} \sum_{i=1}^N (\boldsymbol{e}_i - \sum_{\boldsymbol{e}_k \in n(\boldsymbol{e}_i)} w_k \boldsymbol{e}_k)^2$$

Algorithm

- **1**. Find the weights w_k
- 2. Fix weights w_k and find optimal embeddings e_i solved using a (sparse) eigen-decomposition







the university of edinburgh





Advantages

• globally optimal



Advantages

- globally optimal
- only bottleneck—finding e_i



Advantages

- globally optimal
- only bottleneck—finding e_i
- *not* preserving specific distance





Advantages

- globally optimal
- only bottleneck—finding e_i
- *not* preserving specific distance





Advantages

- globally optimal
- only bottleneck—finding e_i
- *not* preserving specific distance

Issues

reliance on local smoothness





Advantages

- globally optimal
- only bottleneck—finding e_i
- *not* preserving specific distance

Issues

- reliance on local smoothness
- sensitive to noise



he university of edinburgh



Advantages

- globally optimal
- only bottleneck—finding e_i
- *not* preserving specific distance

Issues

- reliance on local smoothness
- sensitive to noise
- no theoretical guarantees about manifold



he university of edinburgh

Project data from X to \mathcal{E} while preserving the *probability distribution over pairwise similarities* between points in the space.



Project data from X to \mathcal{E} while preserving the *probability distribution over pairwise similarities* between points in the space.

$$p(\boldsymbol{x}_{j}|\boldsymbol{x}_{i}) = \frac{\exp(-||\boldsymbol{x}_{i} - \boldsymbol{x}_{j}||^{2}/2\sigma_{i}^{2})}{\sum_{k \neq i} \exp(-||\boldsymbol{x}_{i} - \boldsymbol{x}_{k}||^{2}/2\sigma_{i}^{2})} \quad (i \neq j)$$

$$p(\boldsymbol{x}_{i}|\boldsymbol{x}_{i}) = 0 \qquad \sum_{i} p(\boldsymbol{x}_{j}|\boldsymbol{x}_{i}) = 1$$

$$\mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) = \frac{1}{2N} \left(p(\boldsymbol{x}_{i}|\boldsymbol{x}_{j}) + p(\boldsymbol{x}_{j}|\boldsymbol{x}_{i}) \right)$$

$$\mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{i}) = 0 \qquad \sum_{ij} \mathcal{D}_{X}(\boldsymbol{x}_{i}, \boldsymbol{x}_{j}) = 1$$



Project data from X to \mathcal{E} while preserving the *probability distribution over pairwise similarities* between points in the space.

$$p(\mathbf{x}_{j}|\mathbf{x}_{i}) = \frac{\exp(-||\mathbf{x}_{i} - \mathbf{x}_{j}||^{2}/2\sigma_{i}^{2})}{\sum_{k \neq i} \exp(-||\mathbf{x}_{i} - \mathbf{x}_{k}||^{2}/2\sigma_{i}^{2})} \quad (i \neq j) \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{j}) \coloneqq \operatorname{Student} - t(||\mathbf{e}_{i} - \mathbf{e}_{j}||^{2}, v = 1)$$

$$p(\mathbf{x}_{i}|\mathbf{x}_{i}) = 0 \qquad \sum_{i} p(\mathbf{x}_{j}|\mathbf{x}_{i}) = 1 \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{j}) = \frac{\left(1 + ||\mathbf{e}_{i} - \mathbf{e}_{j}||^{2}\right)^{-1}}{\sum_{k} \sum_{l \neq k} \left(1 + ||\mathbf{e}_{k} - \mathbf{e}_{l}||^{2}\right)^{-1}}$$

$$\mathcal{D}_{\mathcal{X}}(\mathbf{x}_{i}, \mathbf{x}_{j}) = \frac{1}{2N} \left(p(\mathbf{x}_{i}|\mathbf{x}_{j}) + p(\mathbf{x}_{j}|\mathbf{x}_{i}) \right) \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{i}) = 0$$

$$\operatorname{Student} t \text{ allows dissimilar points in } X \text{ to be modelled far away in } \mathcal{E}!$$



Project data from X to \mathcal{E} while preserving the *probability distribution over pairwise similarities* between points in the space.

$$p(\mathbf{x}_{j}|\mathbf{x}_{i}) = \frac{\exp(-||\mathbf{x}_{i} - \mathbf{x}_{j}||^{2}/2\sigma_{i}^{2})}{\sum_{k \neq i} \exp(-||\mathbf{x}_{i} - \mathbf{x}_{k}||^{2}/2\sigma_{i}^{2})} \quad (i \neq j) \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{j}) \coloneqq \operatorname{Student} - t(||\mathbf{e}_{i} - \mathbf{e}_{j}||^{2}, v = 1)$$

$$p(\mathbf{x}_{i}|\mathbf{x}_{i}) = 0 \qquad \sum_{i} p(\mathbf{x}_{j}|\mathbf{x}_{i}) = 1 \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{j}) = \frac{\left(1 + ||\mathbf{e}_{i} - \mathbf{e}_{j}||^{2}\right)^{-1}}{\sum_{k} \sum_{l \neq k} \left(1 + ||\mathbf{e}_{k} - \mathbf{e}_{l}||^{2}\right)^{-1}}$$

$$\mathcal{D}_{\mathcal{X}}(\mathbf{x}_{i}, \mathbf{x}_{j}) = \frac{1}{2N} \left(p(\mathbf{x}_{i}|\mathbf{x}_{j}) + p(\mathbf{x}_{j}|\mathbf{x}_{i}) \right) \qquad \mathcal{D}_{\mathcal{E}}(\mathbf{e}_{i}, \mathbf{e}_{i}) = 0$$

$$Student - t \text{ allows dissimilar points in } X \text{ to be modelled far away in } \mathcal{E}!$$

Objective: $\min \operatorname{KL} \left(\mathcal{D}_{\mathcal{X}}(\boldsymbol{x}_i, \boldsymbol{x}_j) \| \mathcal{D}_{\mathcal{E}}(\boldsymbol{e}_i, \boldsymbol{e}_j) \right)$



t-SNE: Examples





Figures: https://jlmelville.github.io/uwot/umap-examples.html

t-SNE: Examples





t-SNE: Examples





• Different (relatively simple) aspects may be preserved in mapping X to \mathcal{E}



- Different (relatively simple) aspects may be preserved in mapping X to \mathcal{E}
- Trade-off in terms of simplicity of assumption and ease of optimisation



- Different (relatively simple) aspects may be preserved in mapping X to \mathcal{E}
- Trade-off in terms of simplicity of assumption and ease of optimisation
- *t*-SNE quite popular
 - can incorporate both global and local structure (distributional)
 - works on large scale *and* high-dimensional data



- Different (relatively simple) aspects may be preserved in mapping X to \mathcal{E}
- Trade-off in terms of simplicity of assumption and ease of optimisation
- *t*-SNE quite popular
 - can incorporate both global and local structure (distributional)
 - works on large scale *and* high-dimensional data



- Different (relatively simple) aspects may be preserved in mapping X to \mathcal{E}
- Trade-off in terms of simplicity of assumption and ease of optimisation
- *t*-SNE quite popular
 - can incorporate both global and local structure (distributional)
 - works on large scale and high-dimensional data

Issue

None of these methods learn an explicit metric

•
$$\{x_1,\ldots,x_N\} \stackrel{f}{\longrightarrow} \{e_1,\ldots,e_N\}$$
 but $x_t \stackrel{X}{\longrightarrow} e_t$ for unseen x_t

• cannot project an unseen point without redoing the optimisation!



Uniform Manifold Approximation & Projection (UMAP)

Leverage Riemannian geometry and topology to construct a general framework for manifold learning and dimensionality reduction.



Uniform Manifold Approximation & Projection (UMAP)

Leverage Riemannian geometry and topology to construct a general framework for manifold learning and dimensionality reduction.

Overview

0-simplex 1-simplex 2-simplex 3-simplex

• Use simplices as basic building block


Uniform Manifold Approximation & Projection (UMAP)

Leverage Riemannian geometry and topology to construct a general framework for manifold learning and dimensionality reduction.

Overview

- Use simplices as basic building block
- Estimate connectivity and distances between points







Uniform Manifold Approximation & Projection (UMAP)

Leverage Riemannian geometry and topology to construct a general framework for manifold learning and dimensionality reduction.

Overview

- Use simplices as basic building block
- Estimate connectivity and distances between points
- Construct graph that captures topology of manifold!





• Construct faithful topological representation of data



- Construct faithful topological representation of data
- Compute *cross-entropy* between topological structures of X and E in terms of the simplices (building blocks)



- Construct faithful topological representation of data
- Compute *cross-entropy* between topological structures of X and \mathcal{E} in terms of the simplices (building blocks)
- Optimise low-dimensional representation to have minimum cross-entropy



- Construct faithful topological representation of data
- Compute *cross-entropy* between topological structures of X and \mathcal{E} in terms of the simplices (building blocks)
- Optimise low-dimensional representation to have minimum cross-entropy



- Construct faithful topological representation of data
- Compute *cross-entropy* between topological structures of X and E in terms of the simplices (building blocks)
- Optimise low-dimensional representation to have minimum cross-entropy

Advantages

- Can learn a metric, so computing $x_t \xrightarrow{f} e_t$ for unseen x_t is feasible (when labels available)!
- Is fast and scalable
- Can be run unsupervised, supervised, or even weakly supervised!



UMAP: Examples and Comparisons









UMAP: Examples and Comparisons



UMAP: Examples and Comparisons



Figures: https://jlmelville.github.io/uwot/umap-examples.html



• Non-linear dimensionality reduction helps visualise complex data in low dimensions—relies on the manifold hypothesis



- Non-linear dimensionality reduction helps visualise complex data in low dimensions—relies on the manifold hypothesis
- Can explore ways to transform data to run linear versions of algorithms (e.g. kernel PCA)



- Non-linear dimensionality reduction helps visualise complex data in low dimensions—relies on the manifold hypothesis
- Can explore ways to transform data to run linear versions of algorithms (e.g. kernel PCA)
- Most methods exploit the nearest neighbour graph in some form (e.g. MDS, Isomap, LLE, etc.)



- Non-linear dimensionality reduction helps visualise complex data in low dimensions—relies on the manifold hypothesis
- Can explore ways to transform data to run linear versions of algorithms (e.g. kernel PCA)
- Most methods exploit the nearest neighbour graph in some form (e.g. MDS, Isomap, LLE, etc.)
- Data is typically required to be clean (not much noise) and dense ...not too strong a requirement, especially with vision or language



- Non-linear dimensionality reduction helps visualise complex data in low dimensions—relies on the manifold hypothesis
- Can explore ways to transform data to run linear versions of algorithms (e.g. kernel PCA)
- Most methods exploit the nearest neighbour graph in some form (e.g. MDS, Isomap, LLE, etc.)
- Data is typically required to be clean (not much noise) and dense ...not too strong a requirement, especially with vision or language
- Many approaches to choose from—*t*-SNE and UMAP most popular

